

# AutoCollage

Carsten Rother, Lucas Bordeaux, Youssef Hamadi, Andrew Blake  
Microsoft Research Cambridge, UK\*



Figure 1: **AutoCollage automatically creates a collage of representative elements from a set of images.** Novel and desirable properties include: boundaries between images are appropriately positioned; there is little duplication of material; small and meaningless image fragments are avoided; faces are preserved whole; blends may either cut along natural boundaries or be transparent, decided automatically.

## Abstract

The paper defines an automatic procedure for constructing a visually appealing collage from a collection of input images. The aim is that the resulting collage should be representative of the collection, summarising its main themes. It is also assembled largely seamlessly, using graph-cut, Poisson blending of alpha-masks, to hide the joins between input images. This paper makes several new contributions. Firstly, we show how energy terms can be included that: encourage the selection of a representative set of images; that are sensitive to particular object classes; that encourage a spatially

efficient and seamless layout. Secondly the resulting optimization poses a search problem that, on the face of it, is computationally infeasible. Rather than attempt an expensive, integrated optimization procedure, we have developed a sequence of optimization steps, from static ranking of images, through region of interest optimization, optimal packing by constraint satisfaction, and lastly graph-cut alpha-expansion. To illustrate the power of AutoCollage, we have used it to create collages of many home photo sets; we also conducted a user study in which AutoCollage outperformed competitive methods.

\*e-mail: {carrot, lucasb, youssefh, ablake}@microsoft.com

**CR Categories:** I.3.3 [COMPUTER GRAPHICS]: Picture/Image Generation—Display algorithms; I.3.6 [COMPUTER GRAPHICS]: Methodology and Techniques—Interaction techniques; I.4.6 [IMAGE PROCESSING AND COMPUTER VISION]: Segmentation—Pixel classification; partitioning

**Keywords:** Image editing, photomontage, graph cut, energy minimization, constraint satisfaction, Poisson blending

## 1 Introduction

The aim of the paper is to define an automatic procedure for constructing a seamless collage from a collection of input images. There have been various previous studies into related tasks. Pois-

son blending [Perez et al. 2003] can be used to assemble a variety of objects shot against compatible backgrounds onto one seamless background. Digital Photomontage [Agarwala et al. 2004] assembles patches from a batch of repeat shots of a scene, into one composite scene. Graphcut texture [Kwatra et al. 2003], an extension of [Efros and Freeman 2001], synthesises textures by seamlessly joining exemplar patches. In each of the three cases, the pieces to be assembled are broadly compatible, that is to say already approximately matched along the seams, and only adjustment of the seams is needed to render them invisible. Our problem is the one tackled by the Tapestry [Rother et al. 2005] system, in which a set of quite different images has to be composited into a single seamless summarising image. However, Tapestry has certain limitations: it consists of a single optimization step that searches over input images and pixels at every output pixel, and hence is computationally intensive; and it is prone to include small, isolated fragments of the input images. Here we address those limitations by using a multi-stage optimization procedure to tackle complexity, together with explicit region of interest selection. These are the two major improvements of AutoCollage over Tapestry: scalability for large image sets ( $> 50$ ) and robustness. We also include adaptively transparent blending for better hiding of seams, together with object recognition to deal appropriately with sky and faces.

It is also worth mentioning related work on authoring tools for interactive collage systems [Diakopoulos and Essa 2005](and the references therein), where the focus has not been on creating a graphics quality collage, but rather on user interactivity.

**Problem formulation (AutoCollage).** Given a set of input images of arbitrary rectangular shape the aim is to generate a collage of a given rectangular shape, with the following properties.

1. The selected images are representative of the set.
2. From each selected image, one substantial, coherent region of interest (ROI) is extracted.
3. The ROIs should be efficiently packed. Certain objects should be treated with particular respect. In particular, faces should be regarded as preferred material, and should be preserved whole. Sky should be constrained to appear at the top, to avoid sky-like lacunae appearing in the interior of the collage.
4. The transition between images in the collage is subtle, visually smooth and uses transparency where appropriate. The effect we seek is that subjects from the images should be captured and displayed on a background that appears more or less seamless with transitions of input images undetectable.

It may seem that these properties are chosen arbitrarily, however, user feedback (see section 5) suggests that they seem appropriate for this task.

In order to generate an AutoCollage we will express the problem formally as an energy minimization problem, in which each of the desiderata above is represented by an energy term.

The next section introduces the AutoCollage Framework as an energy minimization problem. Section 3 describes our optimization process and includes the main technical contributions. Section 4 explains a new Poisson blending technique that handles transparency. Finally, section 5 displays results of AutoCollage construction, and presents a user study.

## 2 The AutoCollage Framework

The input to AutoCollage is a set of input images  $\mathcal{I} = \{I_1, \dots, I_N\}$ . In order to standardise the input, a pre-processing step is assumed to have been applied, so that each image  $I_n$  is scaled to have unit area, while preserving the aspect ratios of individual images. Then,

building on earlier studies [Rother et al. 2005; Agarwala et al. 2004], on which the proposed framework is based, AutoCollage is viewed as a labelling problem, described using the following notation. The collage is itself an image  $I$ , defined over a domain  $\mathcal{P}$ , and each pixel-location  $p \in \mathcal{P}$  of the collage is to be assigned a label  $L(p)$ , by the AutoCollage algorithm. The labelling  $L = \{L(p), p \in \mathcal{P}\}$  completely specifies the collage, as follows. An individual label has the form  $L(p) = (n, \mathbf{s})$  in which  $I_n \in \mathcal{I}$  is the input image from which the collage pixel  $p$  is taken, and  $\mathbf{s} \in \mathcal{S}$  is the pixel-wise 2D shift of the input image  $n$  with respect to the collage, so that  $I(p) = I_n(p - \mathbf{s})$ . We will often write this compactly as  $I(p) = S(p, L(p))$ , in which  $S(\dots)$  is defined by  $S(p, (n, \mathbf{s})) = I_n(p - \mathbf{s})$  and normalized as  $S(\dots) \in [0, 1] \times [0, 1] \times [0, 1]$ .

The goal of AutoCollage is to find the best labelling  $L \in \mathcal{L}$ , in the space  $\mathcal{L}$  of possible labellings. This is expressed as finding the labelling  $L$  which minimises an energy or cost  $E(L)$ , to be defined in detail later in this section. In the following section, an optimization procedure is defined that searches efficiently in the space of allowed labellings, to obtain a labelling with low energy but, since the algorithm is approximate, not necessarily the global minimum. Note that, by comparison, in Digital Photomontage [Agarwala et al. 2004], all input images were pre-aligned, and therefore each pixel-label consisted of an image index alone, without any shift variable  $\mathbf{s}$ . In AutoCollage, the optimization problem is more complex, because it is necessary to search not only over image indices  $n = 1, \dots, N$ , at each pixel, but also over allowed shifts  $\mathbf{s}$ .

### 2.1 Collage energy

The energy of a labelling  $L$  comprises four terms, as follows:

$$E(L) = E_{\text{rep}}(L) + w_{\text{imp}}E_{\text{imp}}(L) + w_{\text{trans}}E_{\text{trans}}(L) + w_{\text{obj}}E_{\text{obj}}(L) \quad (1)$$

The first term  $E_{\text{rep}}$  tends to select the images from the input image set that are most representative, in two senses: first that chosen images are texturally ‘‘interesting’’ and second that they are mutually distinct so that near duplicates will not be selected. The  $E_{\text{imp}}$  term ensures that a substantial and interesting region of interest (ROI) is selected from each image in  $\mathcal{I}$ . Next,  $E_{\text{trans}}$  is a pairwise term which penalises any transition between images that is not visually appealing. Finally,  $E_{\text{obj}}$  incorporates information on object recognition, and favours placement of objects in reasonable configurations (faces preserved whole, sky at the top, in our implementation). In the remainder of the section, each of these energy terms is defined in detail, together with constraints that must be maintained.

**Representative Image Set** The cost associated with the set  $\mathcal{I}$  of chosen images is of the form  $E_{\text{rep}} = \sum_n E_{\text{rep}}(n)$  where

$$E_{\text{rep}}(n) = -a_n D_r(n) - \min_{m: I_m \in \mathcal{I}} a_n a_m V_r(n, m) \quad (2)$$

and  $a_n$  is an auxiliary, indicator variable, taking the value 1 if the image  $I_n$  is present in the collage and 0 otherwise:

$$a_n = 1 \text{ if } \exists p \in \mathcal{P} \text{ with } L(p) = (n, \mathbf{s}).$$

The unary term  $D_r(n)$  is a measure of the information in image  $n$ . The information measure is defined by

$$D_r(n) = \text{Entropy}(I_n) + w_{\text{face}} \delta(\{\text{Image } n \text{ contains a face}\}) \quad (3)$$

where  $\delta(\pi) = 1$  if predicate  $\pi$  is true, and  $w_{\text{face}}$  weights the influence of an image containing a face, relative to the general textural information in the image. [The histogram used to compute entropy for a given image is constructed in two-dimensional  $a, b$  space from the  $L, a, b$  color system, and discretized into  $16 \times 16$  bins.]

The second term in (2) is expressed in terms of pairwise distances  $V_r(m, n)$  between images, and sums the distances from each

image to its nearest neighbour in the set  $\mathcal{I}$ . As a distance measure  $V_r \in [0, 1]$  we are using normalized chi-squared distance between the color histograms of a pair of images. The histograms are constructed in  $a, b$  space, as above. As well as favouring the most representative images, this energy encourages the use of as many images as possible.

**Importance Cost.** The importance cost consists of a unary term of the form:

$$E_{\text{imp}}(L) = - \sum_p E_{\text{imp}}(p, L(p)). \quad (4)$$

The function  $E_{\text{imp}}(p, L(p)) = G(p, L(p))T(p, L(p))$ , where  $T(p, L(p))$  measures the local entropy, in  $ab$  coordinates, of a  $(32 \times 32)$  pixel region around the pixel  $p$ , and normalised so that local entropy sums to 1 over a given input image. The Gaussian weighting function  $G(\dots)$  favours the centre of the input image from which  $p$  is drawn. Alternatively, instead of  $T$  a more complex model of saliency can be used, as introduced by [Itti et al. 1998].

**Transition Cost.** We use a transition cost similar to those used in the Graphcut texture [Kwatra et al. 2003] and Photomontage [Agarwala et al. 2004] systems. The transition cost is of the form  $E_{\text{trans}} = \sum_{p, q \in N} V_T(p, q, L(p), L(q))$  where  $N$  is the set of all pairs of neighboring (8-neighborhood) pixels. We define the term  $V$  as:

$$V_T(p, q, L(p), L(q)) = \min \left( \frac{\|S(q, L(p)) - S(q, L(q))\|}{\varepsilon + \|S(p, L(p)) - S(q, L(p))\|}, \frac{\|S(p, L(p)) - S(p, L(q))\|}{\varepsilon + \|S(p, L(q)) - S(q, L(q))\|} \right) \quad (5)$$

where intensity function  $S(\dots)$  is as defined above,  $\varepsilon = 0.001$  prevents underflow, and  $\|\cdot\|$  defines the Euclidean norm.

In total,  $E_{\text{trans}}$  measures mismatch across the boundary between two input images. To see this, first observe that  $V_T(p, q, L(p), L(q)) = 0$  unless  $L(p) \neq L(q)$ . Then note that  $V_T(p, q, L(p), L(q))$  is small if there is a strong gradient in one of the input images, since the relevant denominator will then be large. This energy is as in Graph Cut Texture [Kwatra et al. 2003], except that the min operation replaces summation in the original. This is done because, distinctively, adjacent images in this problem are typically taken from rather different scenes, which often do not match. Our choice of energy then acts appropriately in encouraging transition on a high contrast boundary in either scene, in addition to the usual effect of encouraging a good match across the boundary.

**Object Sensitivity** We use state of the art techniques for face detection [Viola and Jones 2001] and general object detection [Shotton et al. 2006] for labelling sky. We would like to exploit this knowledge in such a way that if a face is included, it is included as a whole, and that sky is likely to appear only at the top border of the collage. For faces, as in [Rother et al. 2005], we have the energy term  $E_{\text{obj}} = \sum_{p, q \in N} f(p, q, L(p), L(q))$ , where  $f(p, q, L(p), L(q)) = \infty$  whenever  $L(p) \neq L(q)$  and  $p, q$  are pixels from the same face in either the images of  $L(p)$  or  $L(q)$ , 0 otherwise. For sky rather than defining an explicit energy, we simply label [Shotton et al. 2006] images containing sky and pass this information to the constraint satisfaction engine (see next section) which attempts to position such images only at the top of the collage.

**Parameters** The parameters  $w_{\text{imp}}, w_{\text{trans}}, w_{\text{obj}}, w_{\text{face}}$  have been adjusted by informal testing over 50 sets of home-photographs, where each set contains between 20 – 100 pictures, to achieve reasonably intuitive rankings of the image sets. We take  $w_{\text{imp}} = 10.0, w_{\text{trans}} = 1.0, w_{\text{obj}} = 1.0, w_{\text{face}} = 0.01$ .

**Constraints** The optimization of  $E(L)$  is done under certain constraints, as listed below.

**1. Information bound** Any image  $I_n$  that is present in the labelling, *i.e.* for which  $L(p) = (n, \mathbf{s})$  for some  $\mathbf{s}$  and some  $p \in \mathcal{P}$  must satisfy

$$E_{\text{imp}}(L, n) > T, \quad (6)$$

where  $E_{\text{imp}}(L, n) \in [0, 1]$  is the proportion of local image information  $\sum_p E_{\text{imp}}(p, L(p))$ , as defined above, that is captured in the ROI. In practice we set  $T = 0.9$  — *i.e.* so that at least 90% of the image information is captured. The purpose of this constraint is to guard against the possibility that only a small and unrecognisable fragment of an image may be selected — a problem that plagues the Tapestry system [Rother et al. 2005] — see figure 2. Levying a cost



Figure 2: **Problems with Tapestry.** A collage produced using Tapestry [Rother et al. 2005] includes sky portions in the collage center and other small image fragments. Running AutoCollage on the same image set gives a superior result, see figure 5.

for fragments is quite simply infeasible in the Tapestry framework since it leads to a Markov Random Field with very large cliques, where standard methods such as graph cuts or Belief Propagation are no longer applicable. Here however this is possible thanks to the explicit constraint satisfaction step, which is one of the main innovations of this work — see figure 5.

**2. Uniform shift** A given input image  $I_n$  may appear in the collage with only one unique shift  $\mathbf{s}$ . *i.e.* given two distinct pixels  $p, q \in \mathcal{P} : p \neq q$ , with labels  $L(p) = (n, \mathbf{s}), L(q) = (n, \mathbf{s}')$ , it is required that  $\mathbf{s} = \mathbf{s}'$ . This constraint [Rother et al. 2005] is useful partly for computational efficiency, and partly to ensure that the structure of input images is preserved, without introducing warps.

**3. Connectivity** Each set  $S_n \in \{p \in \mathcal{P} : L(p) = (n, \mathbf{s}), \text{ for some } \mathbf{s}\}$  of collage pixels drawn from image  $n$ , should form a 4-connected region. This cannot in practice be imposed as a hard constraint, but can be encouraged during optimization.

### 3 Energy Minimization

The search space for optimization of the energy  $E(L)$  defined in the previous section, is the entire space of labellings  $L \in \mathcal{L}$ . At each pixel, the input image and its shift must be selected, resulting in a large state-space for graph-cut optimization [Rother et al. 2005]. Here a heuristic but effective approach to optimizing energy  $E(L)$  is adopted in which the various aspects of the labelling are optimized independently and in sequence. First images are ranked statically; then rectangular ROIs are chosen optimally for each image; then a packing problem is solved to assemble and position as many images with highest rank, into the area allowed for the collage, without allowing ROIs to overlap; finally graph-cut optimization fixes pixel identity in areas of overlap of two or more images. Thus the

number of labels for graph-cut is restricted typically to two, three or four, in the overlap areas, exactly as in the Photomontage problem [Agarwala et al. 2004], and thus complexity is dramatically reduced compared with Tapestry [Rother et al. 2005]. This is simply because graph-cut no longer has to optimize over images offsets  $s$ ; those offsets are now determined by optimal packing.

Each of the four optimization steps is described next.

**Image ranking** The ranking step, in the sequence of optimizations, addresses the  $E_{\text{rep}}$  term in the collage energy (1). First images  $I_n$  are relabelled, so that the index  $n$  ranks them according to how representative the subset  $I_1, \dots, I_n$  is. This is straightforward since  $E_{\text{rep}}(n)$  is simply a static rank computed independently in terms of the  $n^{\text{th}}$  image and its predecessors of higher rank. Thus the  $n^{\text{th}}$  image is selected greedily as the one that minimizes

$$-a_n D_{\text{r}}(n) - \min_{m < n} a_n a_m V_{\text{r}}(n, m),$$

adapting the term  $E_{\text{rep}}(n)$  (2). The resulting ranking is then passed to the constraint satisfaction step below.

**ROI optimization** The ROI for each input image  $I_n$  is fixed by minimising, independently for each image, the area of the ROI subject to meeting the information-bound constraint (6), and the constraint that all detected faces are included. This is achieved by constructing a summed area table [Crow 1984] for rapid lookup of the total information  $\sum_{p \in R} E_{\text{imp}}(p, L(p))$  in any rectangular ROI  $R$ . All rectangles are then enumerated, and checked for satisfaction of the constraint, in order to select the one with minimum area. This operation is quadratic in the number of pixels in  $I_n$ , and this is mitigated by subsampling. This is done under the constraint that all detected faces are included. Figure 3 illustrates the effectiveness for this procedure in selecting a ROI.



Figure 3: **ROI selection.** ROIs are determined by selecting the rectangle that optimises  $E_{\text{imp}}$ , and this favours highly textured areas, including all faces, with central positioning.

Note that alternative speed up tricks for the same ROI detection problem have been discussed in [Suh et al. 2005]. Furthermore, they have shown that cropped images, based on the ROI, are more effective for image retrieval compared to using the original images. The only difference to our ROI detection approach is the information measurement where theirs is based on [Itti et al. 1998].

**Constraint satisfaction** Related packing problems have been addressed elsewhere, for example for automatic tiling [Kim and Pellacini 2002]. Here, the packing sub-problem can be stated as follows. We are given a set of selected images and their ROIs, together with the ranking computed above. The goal is to incorporate as many highly ranked images as possible within the width and height of the collage, while respecting the additional constraint that every pixel be covered by some image (though not necessarily covered by some ROI).

The packing problem is a purely combinatorial problem which is reminiscent of applications found in other areas like scheduling [Aggoun and Beldiceanu 1993]. What makes it quite unusual, though, is the simultaneous presence of constraints for *non-overlapping* — no two ROIs should intersect — and *covering* —

every pixel is covered by an images, though not necessarily by a ROI. Being a generalization of well-studied packing problems, it is clear that the problem is NP-hard and that heuristic search is necessary. The general approach is to state the problem as a set of constraints (inequalities, Boolean and linear expressions) between a set of variables. In this problem, the set of variables is

$$\mathcal{V} = \{(x_n, y_n, b_n), n = 1, \dots, N\}, \quad (7)$$

the positions  $(x_n, y_n)$  for each images and a boolean flag  $b_n$  indicating whether the image is to be included or not. Constraints are applied pairwise to images; a typical constraint would be:

$$\text{if } b_n \text{ and } b_m \text{ then } \pi_1 \text{ or } \pi_2, \dots, \quad (8)$$

where a typical proposition is  $\pi_1 = (x_n - x_m > w_m + w_n)$ , in which  $w_m$  and  $w_n$  are respectively the half-widths of the ROIs. Because the relative positions of a ROI pair may be switched, these constraints appear in disjunctive sets — a significant difference from [Kim and Pellacini 2002]. This also puts the problem outside the scope of standard techniques such as Linear Programming. However such problems are amenable to approaches based on constraint programming (CP) [Dechter 2003]. Another hard issue is that constraints are mixed boolean and real, as above, so that effectively constraints must switch between activity and inactivity during optimization. Further object-sensitive constraints can be included — for instance we insist that images with sky appear only at the top of the collage.

To obtain good solutions efficiently, a two-step approach has been used.

**1. Branch and bound** The framework for the first optimization step is a depth-first search which aims at maximising the number of selected images and their quality (Eq. (2)). Constraint propagation [Waltz 1975] is applied to subtrees, from which the subtree may either be pruned, or have its search space reduced. Real variables  $(x_n, y_n)$  are dealt with by coarse discretization with conservative truncation of constraints. The issue of switching the set of active constraints from propagation is dealt with by *reification* [Marriott and Stuckey 1998]. In the branch and bound step, no account is taken of the covering requirement. At this stage we simply solve the problem of packing as many rectangles as possible, within the disjunctive constraints on overlap of ROIs. Even with coarse discretization, the branching factor at a node is large. This is dealt with by randomly selecting a limited number of branches to explore, and by allowing, for each of them, a bounded number of backtracking steps.

**2. Local search** Once branch and bound has terminated, the resulting packing satisfies the non-overlap constraints between ROIs, but in general will not satisfy the covering constraint. At this point, a deterministic local search is applied in order to repair the solution. Perturbations are applied only to  $(x_n, y_n)$ , not to  $b_n$ , so the set of selected images is fixed during this step.

To make sure that a solution which satisfies both the non-overlapping and covering constraints is systematically found, we repeat steps 1) and 2) several times if necessary, and each time relax slightly the constraints (propositions  $\pi_i$  in Eq. 8). In principle, the constraint satisfaction step could generate multiple solutions. After refinement in step 2, these multiple solutions could be evaluated using a bound on the energy function (Eq. 1) or given directly to graph cut optimization. (A bound is needed because strictly the energy function itself is only defined for single coverings of pixels, not multiple coverings as delivered by constraint satisfaction.)

An illustration of the output of the whole constraint satisfaction procedure is given in figure 4, for the collage problem of fig 1.

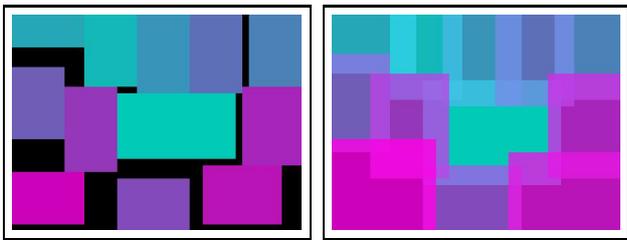


Figure 4: **Output from constraint propagation.** ROIs are shown (left), packed without overlap. Images (right) do overlap, and all pixels are covered.

**Graph cut with alpha expansion** As explained above, graph cut optimization need be applied only to the image variable  $n$  in each pixel-label  $L(p)$  since the shift  $\mathbf{s}$  for each image is now fixed. In practice, up to four values of  $n$  need to be considered at each  $p$  so alpha-expansion [Boykov et al. 2001] is used, exactly as in Digital Photomontage [Agarwala et al. 2004]. Here the objective function to be minimized is that part of the energy  $E$  in (1) that is still “in play”, namely  $w_{\text{imp}}E_{\text{imp}}(L) + w_{\text{trans}}E_{\text{trans}}(L) + w_{\text{obj}}E_{\text{obj}}(L)$ . The first of these terms is unary and the second and third are binary. Since this energy can be shown to be *non-metric*, the truncated schema of alpha-expansion is used [Rother et al. 2005]. At each iteration of alpha-expansion, the 4-connectedness property is encouraged by dilating the optimally expanded set by one pixel.

## 4 $\alpha$ -Poisson Image Blending

The blending task in AutoCollage is to create a seamless transition between input images that are adjacent in the collage. The challenge is to find a general procedure that creates appropriate transitions under differing conditions. Often adjacent images are quite different in colour and texture, so it is not appropriate simply to search for the least obtrusive join [Kwatra et al. 2003], since this will still form a highly visible seam. The choice is either to exploit a natural edge in one of the images, where an edge exists, or else to aim for a soft, transparent blend. One possibility is the extension of Poisson blending [Perez et al. 2003] to include edge-sensitivity [Agarwala et al. 2004]. However this tends to mix colours, rather than achieving the clean transparent effect that we seek here. The solution proposed here is to perform edge-sensitive blending but in the  $\alpha$ -channel rather than in image colour channels.

This is done by computing an alpha mask for each individual input image. In a first step, for a particular image  $I_k$  an overlap area is computed which comprises of all pixels  $p$  where the set of labels  $L(p)$ , which is the same as for the preceding graph-cut optimization, includes label  $I_k$  and at least one other label. Then the following functional minimizes over the overlap area

$$F(\alpha) = \int \|u(\mathbf{r}) - \alpha(\mathbf{r})\|^2 + w(\mathbf{r}) \|\nabla\alpha\|^2 d\mathbf{r}, \quad (9)$$

where  $w(\mathbf{r}) = \lambda + \beta \exp -\frac{1}{2g^2} \max_n \|\nabla I_n\|$  and  $\max_n$  is taken over the images  $I_n$  present in the overlap. Normalising constant  $g^2$  is a mean-square gradient as in [Rother et al. 2004], and we set  $\lambda = 20$ ,  $\beta = 10$ . The function  $u(\mathbf{r})$  takes the value 1 at a pixel  $p$  if the image label, given by graph-cut, is  $I_k$  and 0 otherwise. This selection then biases  $\alpha$  towards the graph-cut solution. Maximization of the functional  $F$  is subject to boundary conditions that  $\alpha = 0, 1$  over the overlap area, and is computed, as usual, by solving a Poisson equation [Perez et al. 2003]. In a final step each image alpha mask is normalized so that at each pixel  $p$  in the output domain the sum of all defined alpha masks is one. The results, clearly visible in fig. 1, is that both sharp abutments and transparent blends are achieved automatically in a collage.

## 5 Results, User Study and Discussion

Results of AutoCollage are shown in figure 1 and 5<sup>1</sup>. Note the features of AutoCollage present in these examples. Boundaries between images are appropriately positioned, avoiding cutting through interesting material. There is little duplication of material and small, visually meaningless image fragments are avoided. Seams between input images switch automatically between cutting along natural boundaries or blending transparently, according to the presence or absence of underlying sharp edges.

Failure modes include the occasional inclusion of sky fragments in the interior, given that sky detection is not infallible (around 83% accuracy in our system). Sometimes texture edges trigger inappropriately sharp transitions in  $\alpha$ . Occasionally face detection fails, allowing an inappropriate cut. A further limitation of the current system is the lack of user interaction. To achieve this we can take further advantage of the constraint satisfaction by including user-specific constraints such as “include one image of this subset”. Obviously, the user should also be able to move, re-size and swap images, and potentially a brush interface, as in [Agarwala et al. 2004], can be used to include explicit image parts in the collage.

In order to answer the questions whether AutoCollage is useful, visually appealing and an improvement over competitive methods, we have conducted a user study<sup>1</sup>. We asked 17 users who were not involved in this work, for a personal dataset (20-50 photos) of an event e.g. a holiday trip. We have created four different summaries for each set: Tapestry [Rother et al. 2005]; AutoCollage with a maximum of 12 images to make it more compatible with Tapestry; PhotoPile which is the collage result after constraint satisfaction only; and a simple grid of 12 most representative pictures. Users were shown a random set of 7 collections, including theirs, and asked to answer the following questions with 1 (definitely no) to 5 (definitely yes). The averaged results are:

- Would you send this summary to a friend?  
AutoCollage (4.6), Tapestry (3.2), PhotoPile (2.2), Grid (1.7)
- Would you like this summary as a screensaver?  
AutoCollage (3.8), Tapestry (2.5), PhotoPile (1.9), Grid (1.25)
- Would you use this summary as front page of the collection?  
AutoCollage (4.5), Tapestry (2.9), PhotoPile (2.6), Grid (2.1)
- Do you find the summary visually appealing?  
AutoCollage (4.5), Tapestry (3.2), PhotoPile (2.3), Grid (1.7)
- Is it a good visual summary of the set of photos?  
AutoCollage (4.6), Tapestry (3.1), Grid (2.8), PhotoPile (2.6)

This shows that AutoCollage is a useful collage tool, superior to the others, for two tasks: Sharing the collage with friends (4.6), and using the collage as a summary of a collection (4.5). AutoCollage was also voted as the visually most appealing. On average users preferred soft blends, similar to the findings in [Diakopoulos and Essa 2005]: “AutoCollage is much more visually appealing, in contrast to PhotoPile which is obviously a stack of photos with hard boundaries” (user quote). The user study gave also insights into the desired properties specified in sec. 1. Firstly, selecting a substantial ROI from each image is important: “I dislike bad details in Tapestry, in particular things like the piece of tree in the middle of a collage”. Secondly, selecting representative image parts is essential: “AutoCollage is definitely better than Tapestry by having a mix of people and landscape; Tapestry shows sometimes only faces.” Finally, object (sky) recognition is essential: “Having a piece of sky

<sup>1</sup>Further results, and details of the user study, are available on <http://research.microsoft.com/vision/cambridge/i3l/AutoCollage/default.htm>



Figure 5: **Further example of AutoCollage.** Note again, the desirable and novel properties: hidden boundaries, little duplication, freedom from fragments, selective transparency.

in the center of a collage looks funny, I thought it's snow". Further user quotes: "Seeing an AutoCollage of your own photos is a surprisingly emotive experience."; "AutoCollage is a clear win on average; it is better proportioned than the others"; "AutoCollage is sometimes much better than Tapestry, however, AutoCollage is never worse than Tapestry. I'd also like to include a specific image".

The improvement of AutoCollage compared to Tapestry is also prominent in terms of runtime. To create collages in fig. 5 and 2, from a set of 14 input images, AutoCollage is 16 times faster in the core packing stage: 0.13sec (AutoCollage) vs. 2.07sec (Tapestry, with only 0.001% of all possible expansion moves).

In future work it is planned to elaborate the object sensitivity. We will also address the issue of user interactivity which depends, however, heavily on the user task at hand. Another interesting possibility is to allow multiple solutions from constraint satisfaction, and then evaluate each using a suitable bound on the energy.

## References

- AGARWALA, A., DONTCHEVA, M., AGRAWALA, M., DRUCKER, S., COLBURN, A., CURLESS, B., SALESIN, D., AND COHEN, M. 2004. Interactive digital photomontage. *ACM Trans. Graph.* 23, 3, 294–302.
- AGGOUN, A., AND BELDICEANU, N. 1993. Extending CHIP in order to solve complex scheduling and placement problems. *Mathematical Computer Modelling* 17, 7, 57–73.
- BOYKOV, Y., VEKSLER, O., AND ZABIH, R. 2001. Fast approximate energy minimization via graph cuts. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 23, 11.
- CROW, F. 1984. Summed area tables for texture mapping. In *Proc. ACM Siggraph*, ACM, 207–212.
- DECHTER, R. 2003. *Constraint Processing*. Morgan Kaufmann.
- DIAKOPOULOS, N., AND ESSA, I. 2005. Mediating photo collage authoring. In *UIST*, CD-ROM.
- EFROS, A. A., AND FREEMAN, W. T. 2001. Image quilting for texture synthesis and transfer. *Proc. ACM Siggraph*.
- ITTI, L., KOCH, C., AND NIEBUR, E. 1998. A model of saliency based visual attention for rapid scene analysis. *IEEE Trans. on Pattern Analysis and Machine Intelligence* 20, 11.
- KIM, J., AND PELLACINI, F. 2002. Jigsaw image mosaics. In *Proc. ACM Siggraph*, ACM, 657–664.
- KWATRA, V., SCHODL, A., ESSA, I., TURK, G., AND BOBICK, A. 2003. Graphcut textures: image and video synthesis using graph cuts. *ACM Trans. Graph.* 22, 3, 277–286.
- MARRIOTT, K., AND STUCKEY, P. 1998. *Programming with Constraints*. The MIT Press.
- PEREZ, P., GANGNET, M., AND BLAKE, A. 2003. Poisson image editing. *ACM Trans. Graph.* 22, 3, 313–318.
- ROTHER, C., KOLMOGOROV, V., AND BLAKE, A. 2004. Grabcut: Interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.* 23, 3, 309–314.
- ROTHER, C., KUMAR, S., KOLMOGOROV, V., AND BLAKE, A. 2005. Digital tapestry. In *Proc. Conf. Comp. Vision and Pattern Recog.*
- SHOTTON, J., WINN, J., ROTHER, C., AND CRIMINISI, A. 2006. Textonboost: Joint appearance, shape and context modelling for multi-class object recognition and segmentation. In *Europ. Conf. Comp. Vision*.
- SUH, B., LING, H., BEDERSON, B. B., AND JACOBS, D. W. 2005. Automatic thumbnail cropping and its effectiveness. In *UIST*, CD-ROM.
- VIOLA, P., AND JONES, M. 2001. Rapid object detection using a boosted cascade of simple features. In *Proc. Conf. Comp. Vision and Pat. Recog.*
- WALTZ, D. 1975. Understanding line drawings of scenes with shadows. In *The Psychology of Vision*, W. P.H., Ed. McGraw-Hill, New York.